# Double precision is not needed for many-body calculations: New conventional wisdom

Pavel Pokhilko[a], Evgeny Epifanovsky[b], and Anna I. Krylov[a]

[a] Department of Chemistry, University of Southern California, Los Angeles, California 90089-0482

[b] Q-Chem Inc., 6601 Owens Drive, Suite 105 Pleasanton, California 94588

April 2, 2018

**Abstract**

Using single precision floating point representation reduces the size of data and computation time by a factor of two relative to double precision conventionally used in electronic structure programs. For large-scale calculations, such as those encountered in many-body theories, reduced memory footprint alleviates memory and input/output bottlenecks. Reduced size of data can lead to additional gains due to improved parallel performance on CPUs and various accelerators. However, using single precision can potentially reduce the accuracy of computed observables. Here we report an implementation of coupled-cluster and equation-of-motion coupled-cluster methods with single and double excitations in single precision. We consider both standard implementation and one using Cholesky decomposition or resolution-of-the-identity of electron-repulsion integrals. Numerical tests illustrate that when single precision is used in correlated calculations, the loss of accuracy is insignificant and pure single-precision implementation can be used for computing energies, analytic gradients, excited states, and molecular properties. In addition to pure single-precision calculations, our implementation allows one to follow a single-precision calculation by clean-up iterations, fully recovering double-precision results while retaining significant savings.

## 1 Introduction

Quantum chemistry is one of the most demanding fields in terms of computational resources. Standard formulations of many-body theories result in large amount of data (wave-function parameters) and in steep computational scaling. For example, storage and floating point operations requirements of the coupled-cluster method with single and double substitutions (CCSD) scale as $N^4$ and $N^6$ with the system size, respectively.[1] This steep scaling limits the applicability of these highly reliable methods. Large memory footprint, inherent to correlated theories, also creates a hurdle for efficient parallelization and utilization of accelerators (such as graphic processing units, GPUs). In this paper, we present a production-level implementation of CCSD and EOM-CCSD (equation-of-motion CCSD) methods[2–8] and investigate the impact of using reduced precision on computational efficiency and the accuracy of the results.

The standard of accuracy for quantum chemical calculations of thermochemical and chemical kinetics data is 1 kcal/mol (which is equal 4.2 kJ/mol or $1.593 \cdot 10^{-3}$ hartree).[9,10] A typical

desired accuracy for excitation or ionization energies is 0.01-0.1 eV. The errors in the calculations arise due to the intrinsic errors of the methods due to approximations in many-electron and one-electron bases[1] as well as due to the finite convergence thresholds for the iterative algorithms. Typically, convergence thresholds are much tighter than the methods' error bars.

The IEEE 754 standard[11,12] defines single and double precision floating-point arithmetic in computers. The numbers in this format are represented in the scientific notation

$$(-1)^s b_0.b_1 b_2 b_3 \ldots b_{p-1} 2^E, \tag{1}$$

where $s$ is the sign bit, $p$ is the precision of significand, $E$ is the exponent, and bits $b_i$ can take values 0 or 1. Although the second revision of the standard[12] generalizes base, or *radix*, here we consider only binary formats. The number is called *normal* if $b_0$ is 1 and *subnormal* otherwise. Subnormal numbers fill the gap between the smallest positive normal number and zero. The attributes of single- and double-precision floating-point numbers are summarized in Table 1.

Table 1: Summary of single and double floating-point IEEE 754 standard.

|  | Single | Double |
|---|---|---|
| Total size, bits | 32 | 64 |
| Exponent size, bits | 8 | 11 |
| Exponent bias | 127 | 1023 |
| Sign, bits | 1 | 1 |
| Significant (explicit), bits | 23 | 52 |
| Decimal precision | $\log_{10}(2^{24}) \approx 7$ | $\log_{10}(2^{53}) \approx 16$ |
| Smallest normal number | $2^{-126}$ | $2^{-1022}$ |

Typical implementations of most *ab initio* methods use double precision arithmetic. The single precision format uses half the number of bits of double precision, thereby allowing to store twice as many values. Another benefit is proportionally faster memory access and disk input/output (in terms of the number of elements per second). This is important in practice because memory speed improves at a slower rate than CPUs. Furthermore, CPU caches can accommodate twice as many floating point numbers potentially leading to less frequent cache misses. In terms of computational time, single precision gives a twofold speedup on CPUs for most modern architectures and the gains on GPUs can be much larger. Thus, single-precision implementation of quantum-chemistry methods can extend the scope of systems amenable to these treatments, decrease time-to-solution, and reduce energy footprint. Using too much energy and power per calculation is recognized as one of the biggest challenges in high performance computing and using reduced precision or even entirely different representation of real numbers have been advocated.[13]

Although double precision is *de facto* a standard in quantum chemistry and other scientific calculations, many algorithms can be re-designed to work in mixed precision, as was recently done in such diverse areas as lattice quantum chromodynamics,[14] molecular dynamics,[15] and general linear algebra algorithms with iterative refinements.[16] In quantum chemistry mixed-precision algorithms have been explored in the context of integral calculations[17,18] within Hartree-Fock (HF) and density functional theory (DFT).[19–22] The main conclusion from these studies was that pure single precision is not sufficient for integral and HF/DFT calculations and one should only deploy it in a mixed-precision fashion, i.e., such that some operations are

performed in single precision and some — in double precision. The utility of single precision has not yet been thoroughly investigated in post-HF calculations. Previous studies[23–25] focused on non-iterative methods, such as MP2 and triples correction for CCSD. Single-precision MP2 was tested within resolution-of-the-identity (RI)[23] and Cholesky decomposition (CD)[24] schemes. These studies have shown that commonly used RI bases and CD thresholds ($10^{-2}$–$10^{-3}$) yield much larger errors in total energies than errors due to using single-precision arithmetics. Numerical analysis of the (T) correction for CCSD[25] has shown that this calculation is stable with respect to numerical noise, justifying single-precision implementation. Despite these encouraging findings, the extent of applicability of single precision in many-body theories is not fully understood. There are several open questions:

1. Does numerical error accumulate in iterative procedures such as those used to solve CCSD and EOM equations?

2. What is the impact of using single precision on molecular properties and excited states?

3. Can one reliably compute analytic nuclear gradients and optimize structures within pure single precision?

As a standard practice in quantum chemistry, typical numerical convergence thresholds are tight enough to not affect the resulting accuracy. For example, Q-Chem's default CCSD convergence criteria are $10^{-6}$ hartree for energies and $10^{-4}$ for amplitudes;[26] Molpro uses $10^{-6}$ hartree for energies and $10^{-5}$ for amplitudes.[27] Some packages use a single threshold, for example, GAMESS uses $10^{-7}$ for amplitudes[28] and ORCA uses $10^{-5}$–$10^{-6}$ hartree for energies. Thus, it appears that 7 decimal digits is sufficient for correlation energy (we note that $10^{-7}$ hartree is three orders of magnitude tighter that chemical accuracy).

In this paper, we describe a general implementation of CC/EOM-CC methods that allows users to perform calculations in either single or double precision. We implemented both the standard variant and one using Cholesky decomposition (CD) or resolution-of-the-identity (RI) of electron-repulsion integrals. In addition, our implementation allows one to follow up a single-precision calculation with a clean-up step in which the full double-precision accuracy can be recovered. The code is based on the libtensor[29] and libxm[30] libraries for many-body electronic structure calculations. The production-level code is implemented in the Q-Chem electronic structure package.[31, 32]

# 2    Algorithms and implementation details

Libtensor[29] was developed to provide a high-level interface for tensor operations and to deliver efficient performance. The library has been used to implement a large number of CC,[2–4] EOM-CC,[6–8,33] and ADC (algebraic diagrammatic construction)[34] methods for calculating energies, properties, and nuclear gradients in Q-Chem. Libtensor supports tensor symmetries, block-sparcity, several contraction algorithms, different BLAS implementations, and several computational backends.[30,35] We extended the libtensor library to incorporate single-precision operations by generalizing all tensor operations for the general element type. Fig. 1 gives an example of the code. The code is available in the original repository on github. The numerical tests presented in this paper were performed with a developer version of Q-Chem and the

```
template<size_t N, typename T>
class bto_copy :
    public additive_gen_bto<N, typename bto_traits<T>::bti_traits>,
    public noncopyable {
// ....
// Declaration of the templated class
}

template<size_t N>
using btod_copy = bto_copy<N, double>;
```

Figure 1: Example of the change in the interface. In this code a block-tensor operation over double type ("btod") is generalized to a block-tensor operation over the template type.

modified libtensor, compiled by the GCC-6.4.0 compiler with the '-O3' optimization flag and linked to the Intel MKL library (2017.0 version).

Fig. 2 shows an overview of the CCSD and EOM-CCSD workflows. In our implementation, we compute all integrals and solve SCF equations in double precision; the integral transformation step is also performed in double precision. We then convert the tensors (integrals) to single precision and perform all tensor operations (contractions, additions, etc.) in the CCSD/EOM-CCSD calculations using single precision. In the RI/CD variants, we also employ the single precision algorithm at the CCSD step. Once the CCSD equations converge in single precision, the procedure can switch to double precision to perform clean-up iterations to recover the double-precision result. We follow the same strategy for $\Lambda$-amplitudes. We also implemented single-precision calculations of various density matrices needed for property and nuclear gradient calculations.

The EOM workflow entails solving the CCSD equations and evaluating the similarity-transformed Hamiltonian ($\bar{H}$) intermediates followed by the computation of EOM energies and amplitudes by iterative diagonalization of the similarity-transformed Hamiltonian using the Davidson algorithm. We extended the routines[36] that compute the intermediates and $\sigma$-vectors (products of the similarity-transformed Hamiltonian acting on EOM trial states) for EOM-IP/EA/SF/EE-CCSD to support calculations in single precision.

It is expected that using single precision provides a speedup factor of two on CPUs (not taking into account additional speedup due to reduced IO). If the single-precision calculation is followed up by clean-up steps to recover double-precision accuracy, the theoretical estimate for the overall speedup is estimated as:

$$\text{Theor. speedup} = \frac{N_{\text{dp}}}{0.5 \cdot N_{\text{sp}} + N_{\text{cleanup}}}, \tag{2}$$

where $N_{\text{dp}}$ and $N_{\text{sp}}$ are the number of iterations in double and single precisions, $N_{\text{cleanup}}$ is the number of cleanup iterations in double precision.
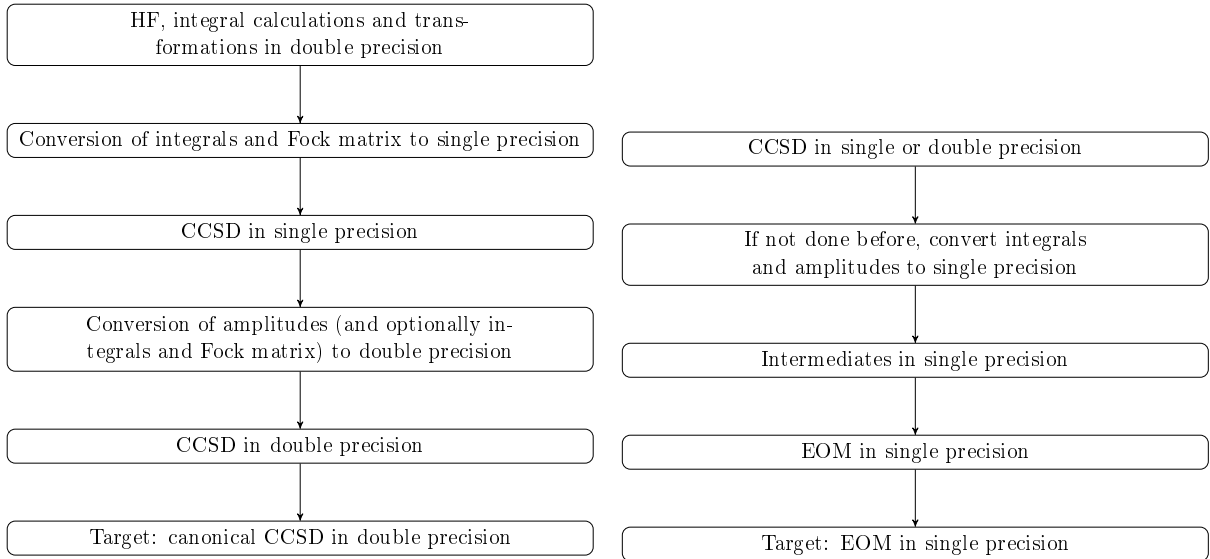
Figure 2: Left: Implemented CCSD algorithm. Cleanup step in double precision is optional. Right: EOM algorithm.

Table 2: Convergence thresholds for CCSD calculations. Convergence for $\Lambda$ equations is the same as for the $T$ amplitudes.

| System | $E$, sp, hartree | $T$, sp | $E$, dp, hartree | $T$, dp |
|---|---|---|---|---|
| Water clusters | $10^{-5}$ | $10^{-3}$ | $10^{-6}$ | $10^{-4}$ |
| Uracil | $10^{-6}$ | $10^{-4}$ | $10^{-6}$ | $10^{-4}$ |
| ATT | $10^{-6}$ | $10^{-4}$ | $10^{-6}$ | $10^{-4}$ |
| G2 set, pure sp or dp | $10^{-6}$ | $10^{-4}$ | $10^{-6}$ | $10^{-4}$ |
| G2 set, sp with cleanup | $10^{-5}$ | $10^{-3}$ | $10^{-6}$ | $10^{-4}$ |
| G2 set, sp with energy in dp | $10^{-5}$ | $10^{-3}$ | — | — |
| $C_6H_5N$ | $10^{-6}$ | $10^{-4}$ | $10^{-6}$ | $10^{-4}$ |
| benzene | $10^{-10}$ | $10^{-9}$ | $10^{-7}$ | $10^{-5}$ |

# 3   Computational details

Benchmark calculations were performed on 4x8 Intel Xeon E5-4640 using 4 threads. The benchmark set consists of diverse types of electronic structure, including water clusters of increasing size, the uracil molecule, a nucleobase trimer (ATT, adenine-thymine-thymine), an aromatic diradical ($C_6H_5N$), the benzene molecule, and the G2 set[37] containing 148 molecules. We tested basis sets of the double-, triple-, and quadruple-zeta quality, with and without diffuse functions. All Cartesian geometries are given in Supplementary Information (SI).

Convergence thresholds for CCSD equations are summarized in Table 2. Since single precision provides ~7 decimal digits of precision, the tightest threshold of convergence for energy in single precision is $10^{-7}$ hartree (assuming correlation energies of 1 hartree). By numerical experimentation we found that single precision iterations converge smoothly with energy threshold of $10^{-5}$–$10^{-6}$ hartree; therefore, we used this threshold for the single precision part of the calculation in most cases. In the double-precision calculations, we used default convergence except for the benzene benchmark. In properties calculations, we used the same convergence

criteria for $T$ and $\Lambda$ amplitudes. Unrelaxed one-particle density matrices used in dipole moment calculations were computed in double and single precision from respective double and single precision intermediates and amplitudes. In all EOM calculations a $10^{-5}$ convergence threshold for EOM amplitudes[38] was used in the Davidson procedure. For geometry optimization and finite-difference frequency analysis, much tighter convergence criteria were used: $10^{-10}$ for energies, $10^{-9}$ for $T$ and $\Lambda$ amplitudes for double precision; $10^{-7}$ for energies, $10^{-5}$ for $T$ and $\Lambda$ amplitudes for single precision. The criteria of convergence in geometry optimization were the same for the double and single precision calculations: $2 \cdot 10^{-5}$ a.u. for the maximum component of the gradient $5 \cdot 10^{-5}$ a.u. for the maximum atomic displacement, and $1 \cdot 10^{-7}$ a.u. for energy.

# 4 Results and discussion

## 4.1 Accuracy of ground-state energies and properties

We found that the CC amplitude convergence rate is not affected by using single precision and that the number of CCSD iterations in calculations with single-precision CCSD followed by a clean-up step is the same as in the reference double-precision calculations. A typical output of single-precision calculation with the clean-up step is shown in SI (Fig. S1).

Table 3 shows the results for water clusters. Pure single precision calculations (without the clean-up step) do not introduce significant numerical errors in total energies, i.e. the typical difference between single and double precision energies is only several J/mol, which is three orders of magnitude below chemical accuracy. Moreover, the double-precision numerical accuracy is fully recovered when the cleanup iteration is performed. The single-precision calculation is twice as fast than the double-precision one. For the calculation with the clean-up step, the speedup quickly approaches the theoretical maximum (about 1.5, for these parameters) with the increasing number of water molecules (see Fig. 3). We observed similar performance and accuracy for RI/CD-CCSD calculations of these clusters.
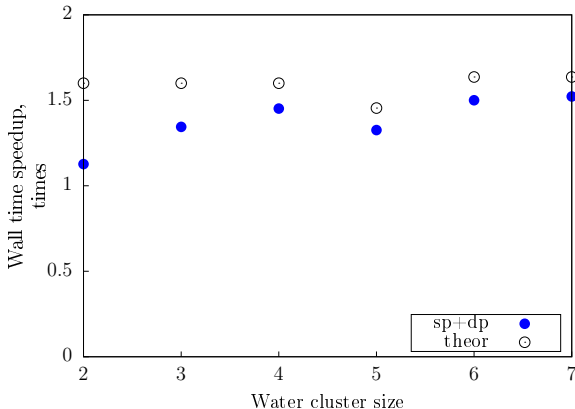


Figure 3: Wall time speedup for SP scheme with cleanup for the CCSD/cc-pVDZ energy calculations of water clusters. Frozen core is used. Theoretical estimate is given by Eq. (2).

To test whether numerical errors increase with system size, we compared single and double precision CCSD for an adenine-thymine-thymine system (ATT) and found that the difference

between single and double precision total energies is 15 J/mol (again, full double precision accuracy can be recovered with a single clean-up step).

We investigated basis set effects on the error due to single precision by using uracil. Table 5 presents total energies and dipole moments. The results show that the errors in energies are negligible and that the differences in dipole moments computed in single and double precision are less than the number of digits printed in the output. Increasing the basis set from cc-pVDZ to cc-pVTZ and aug-cc-pVTZ does not increase the errors.

To confirm that these observations hold for other systems, we compared the CCSD total energies computed in single and double precision for the G2 set[37] using the the 6-31G(d), cc-pVDZ, cc-pVTZ, and cc-pVQZ basis sets. In addition, we compared single-precision results with those obtained with a cleanup step and with a calculation in which energy is computed in double precision using double-precision integrals and single-precision amplitudes. The results are summarized in Table 6. As one can see, the MAD in single-precision calculations is less than 4 J/mol. The MADs, STDs, and maximum errors in cc-pVQZ basis are only slightly larger than in cc-pVDZ.

In order to meaningfully compare the differences between single and double precision, one should keep in mind that in these calculations the energy threshold for CCSD convergence was $10^{-6}$ hartree $\approx 2.6$ J/mol. Within this threshold, the error of the calculations with the cleanup step is the same as in the single-precision calculation. Interestingly, calculation of double-precision energies from single-precision amplitudes gives larger errors, which can be explained by the convergence criteria used: the amplitudes are underoptimized in comparison with the amplitudes, used for other calculations (shown in Table 2).

Table 3: CCSD/cc-pVDZ total energies of water clusters. Differences between single- and double-precision energies are shown in the last column.

| Cluster size | $E_{sp}$, a.u. | $E_{sp+dp}$, a.u. | $E_{dp}$, a.u. | $\Delta_{dp,sp}$, J/mol |
|---|---|---|---|---|
| 2 | -152.48710832 | -152.48710752 | -152.48710767 | 1.7 |
| 3 | -228.74911206 | -228.74911017 | -228.74911059 | 3.9 |
| 4 | -305.00868540 | -305.00868416 | -305.00868472 | 1.8 |
| 5 | -381.26858576 | -381.26858599 | -381.26858585 | -0.2 |
| 6 | -457.51949116 | -457.51949130 | -457.51949145 | -0.8 |
| 7 | -533.78980571 | -533.78980623 | -533.78980640 | -1.8 |

Table 4: CD-CCSD/cc-pVDZ total energies of water clusters. Cholesky threshold of $10^{-3}$ was used. Differences between single and double precision energies are shown in the last column.

| Cluster size | $E_{sp}$, a.u. | $E_{sp+dp}$, a.u. | $E_{dp}$, a.u. | $\Delta_{dp,sp}$, J/mol |
|---|---|---|---|---|
| 2 | -152.48713650 | -152.48713571 | -152.48713586 | 1.7 |
| 3 | -228.74916915 | -228.74916737 | -228.74916779 | 3.6 |
| 4 | -305.00875903 | -305.00875662 | -305.00875717 | 4.9 |
| 5 | -381.26864546 | -381.26864554 | -381.26864540 | 0.2 |
| 6 | -457.51960841 | -457.51960843 | -457.51960858 | -0.4 |
| 7 | -533.78994271 | -533.78994226 | -533.78994243 | 0.7 |

Table 5: CCSD total energies and dipole moments of uracil in various basis sets.

| Basis set | $E_{\text{sp}}$, a.u. | $E_{\text{dp}}$, a.u. | $\Delta_{dp,sp}$, J/mol | Dipole, T,L,1DM in sp | Dipole, T,L,1DM in dp |
|---|---|---|---|---|---|
| cc-pVDZ | -413.72139806 | -413.72139812 | -0.2 | 1.550917 | 1.550917 |
| cc-pVTZ | -414.10170402 | -414.10170403 | 0.0 | 1.643366 | 1.643366 |
| aug-cc-pVTZ | -414.13153216 | -414.13153108 | 2.8 | 1.689405 | 1.689405 |

Table 6: Mean average deviation (MAD) and standard deviation (STD), J/mol, from reference double-precision CCSD total energies for the G2 set.

| Basis | sp | | sp with dp energy | | sp with cleanup | |
|---|---|---|---|---|---|---|
| | MAD | STD | MAD | STD | MAD | STD |
| 6-31G(d) | 0.12 | 0.2 | 3.7 | 6.6 | 1.5 | 1.4 |
| cc-pVDZ | 0.15 | 0.2 | 3.6 | 6.7 | 1.4 | 1.4 |
| cc-pVTZ | 0.3 | 0.5 | 3.4 | 5.3 | 1.5 | 1.4 |
| cc-pVQZ | 0.68 | 0.9 | 3.7 | 5.0 | 1.6 | 1.5 |

## 4.2 Accuracy of target-state energies in EOM-CCSD

To test the accuracy of single-precision calculations of excited states, we carried out EOM-EE-CCSD calculations for uracil and EOM-SF-CCSD calculations for $C_6H_5N$. We considered a range of different states: singlets and triplets, states with different symmetry (symmetric and antisymmetric with respect to the symmetry plane), closed-shell and open-shell types, valence and Rydberg states. In all cases, the difference between total energies computed in double and single precision does not exceed 3 J/mol ($3 \cdot 10^{-5}$ eV), which is much smaller than the intrinsic error bars of EOM-CCSD (0.1-0.3 eV). This result is comparable to the differences in the CCSD total energies. The respective excitation energies are the same within at least 4 decimal figures.

## 4.3 Accuracy of gradient evaluation in single precision

To test how the precision affects the accuracy of the gradient calculation and the resulting optimized structures, we optimized the benzene molecule at the CCSD/cc-pVDZ level of theory using tight convergence criteria. In these calculations, density matrices were computed in the respective precisions; orbital response was computed in double precision. Geometry optimization was performed by gradient optimization (gradient was evaluated analytically). Starting from the same initial geometry (MP2/cc-pVDZ optimized structure), both optimizations converged in 6 iterations (single precision run converged by gradient and displacement). The resulting geometries are nearly identical, with mean absolute error (computed for only non-zero Cartesian coordinates) of $1.9 \cdot 10^{-9}$ Å.

## 4.4 Accuracy of finite-difference frequencies

To investigate whether finite-difference calculations of frequencies (using analytic gradients) can be carried out in single precision, we computed the frequencies and normal modes for benzene at the respective optimized geometries with the same convergence criteria as used for geometry

Table 7: EOM-EE-CCSD total energies (in a.u.) of excited singlet and triplet states of uracil in various basis sets. In each cell the first number is obtained from double-precision calculation and the second number is obtained from single-precision EOM calculation. The third number is the difference between double- and single-precision energies. CCSD equations were solved in double precision in all cases.

| Singlets Basis set | 1A′ | 2A′ | 1A″ | 2A″ |
|---|---|---|---|---|
| cc-pVDZ | -413.50624322 | -413.46584035 | -413.52849406 | -413.47629005 |
| | -413.50624315 | -413.46584026 | -413.52849409 | -413.47628997 |
| | $-7 \cdot 10^{-8}$ | $-9 \cdot 10^{-8}$ | $3 \cdot 10^{-8}$ | $-8 \cdot 10^{-8}$ |
| cc-pVTZ | -413.89150482 | -413.84840838 | -413.90876266 | -413.85792039 |
| | -413.89150498 | -413.84840792 | -413.90876251 | -413.85792058 |
| | $1.6 \cdot 10^{-7}$ | $-4.6 \cdot 10^{-7}$ | $-1.5 \cdot 10^{-7}$ | $1.9 \cdot 10^{-7}$ |
| aug-cc-pVTZ | -413.92599740 | -413.88373275 | -413.94000534 | -413.90496469 |
| | -413.92599743 | -413.88373272 | -413.94000543 | -413.90496414 |
| | $3 \cdot 10^{-8}$ | $-3 \cdot 10^{-8}$ | $9 \cdot 10^{-8}$ | $-5.5 \cdot 10^{-7}$ |
| Triplets Basis set | 1A′ | 2A′ | 1A″ | 2A″ |
| cc-pVDZ | -413.57714822 | -413.51667861 | -413.53936723 | -413.48574474 |
| | -413.57714817 | -413.51667858 | -413.53936718 | -413.48574485 |
| | $-5 \cdot 10^{-8}$ | $-3 \cdot 10^{-8}$ | $-5 \cdot 10^{-8}$ | $1.1 \cdot 10^{-7}$ |
| cc-pVTZ | -413.95932894 | -413.89779284 | -413.91902233 | -413.86657196 |
| | -413.95932869 | -413.89779281 | -413.91902181 | -413.86657206 |
| | $-2.5 \cdot 10^{-7}$ | $-3 \cdot 10^{-8}$ | $-5.2 \cdot 10^{-7}$ | $1.0 \cdot 10^{-7}$ |
| aug-cc-pVTZ | -413.99008713 | -413.92935022 | -413.94969790 | -413.90754309 |
| | -413.99008666 | -413.92935024 | -413.94969820 | -413.90754252 |
| | $-4.7 \cdot 10^{-7}$ | $2 \cdot 10^{-8}$ | $3.0 \cdot 10^{-7}$ | $-5.7 \cdot 10^{-7}$ |

Table 8: EOM-SF-CCSD total energies (in a.u.) of several electronic states of $C_6H_5N$ in various basis sets. Symmetry labels refer to state symmetries. In each cell the first number is from double-precision calculation, the second number is from single-precision EOM calculation, and the third number is the difference between double- and single-precision energies. CCSD equations were solved in double precision in all cases.

| Basis set | 1A$_2$ | 2A$_2$ | 1A$_1$ | 2A$_1$ |
|---|---|---|---|---|
| cc-pVDZ | -285.46570954 | -285.42658952 | -285.40655223 | -285.37250140 |
| | -285.46570956 | -285.42658954 | -285.40655224 | -285.37250139 |
| | $2 \cdot 10^{-8}$ | $2 \cdot 10^{-8}$ | $1 \cdot 10^{-8}$ | $-1 \cdot 10^{-9}$ |
| cc-pVTZ | -285.79837285 | -285.76114321 | -285.74329035 | -285.70728705 |
| | -285.79837303 | -285.76114343 | -285.74328950 | -285.70728629 |
| | $1.8 \cdot 10^{-7}$ | $2.2 \cdot 10^{-7}$ | $-8.5 \cdot 10^{-7}$ | $-7.6 \cdot 10^{-7}$ |

optimizations with the step size of 0.001 Å. The resulting frequencies are very different: real and imaginary frequencies of ∼15000 cm$^{-1}$ occur along with several imaginary frequencies. Thus, not surprisingly, single precision is not sufficient for finite-difference calculations. How-

ever, finite-difference calculations using single-precision calculation with the double-precision cleanup step fully recovers double precision frequencies, while affording $\sim \times 1.3$ speedup of the calculation.

# 5    Conclusion

In this contribution, we report single- and mixed-precision implementation of the CCSD and EOM-CCSD energies, analytic gradients, and properties. Using single precision results in reduced memory footprint, considerable computation speed-up, and reduced energy-to-solution. That is, single precision calculations use half the memory (or disk) space and produce a speedup factor of 2 on CPUs. The main focus of this paper was on assessing the impact on accuracy of the resulting energies and properties, in particular in iterative schemes. The results indicate that the rate of error accumulation in single precision is sufficiently slow, and overall single precision introduces negligible errors for total energies, excitation energies, forces, and properties. Moreover, single precision can be used in geometry optimizations with analytical gradients.

If tight convergence is desired (e.g. in finite-difference calculations), single precision can be used to speedup iterations in the beginning, converging to the single-precision result first and continuing in double precision. In most cases the total number of iterations is not affected and the speedup is close to the theoretical estimate.

# 6    Acknowledgment

Supporting information is available: relevant Cartesian geometries, optimized structures, and frequencies and normal modes.

# References

[1] Helgaker, T.; Jørgensen, P.; Olsen, J. *Molecular electronic structure theory;* Wiley & Sons, 2000.

[2] Bartlett, R.J.; Shavitt, I. *Many-Body Methods in Chemistry and Physics: MBPT and Coupled-Cluster Theory;* Cambridge University Press, 2009.

[3] Bartlett, R.J. The coupled-cluster revolution *Mol. Phys.* **2010**, *108*, 2905–2920.

[4] Bartlett, R.J. In *Theory and Applications of computational chemistry*; Dykstra, C., Frenking, G., Scuseria, G., Eds.; Elsevier, 2005.

[5] Bartlett, R.J. To multireference or not to multireference: That is the question? *Int. J. Mol. Sci.* **2002**, *3*, 579–603.

[6] Krylov, A. I. Equation-of-motion coupled-cluster methods for open-shell and electronically excited species: The hitchhiker's guide to Fock space *Annu. Rev. Phys. Chem.* **2008**, *59*, 433–462.

[7] Sneskov, K.; Christiansen, O. Excited state coupled cluster methods *WIREs Comput. Mol. Sci.* **2012**, *2*, 566–584.

[8] Bartlett, R.J. Coupled-cluster theory and its equation-of-motion extensions *WIREs Comput. Mol. Sci.* **2012**, *2*, 126–138.

[9] Pople, J.A. Nobel lecture: Quantum chemical models *Rev. Mod. Phys.* **1999**, *71*, 1267–1274.

[10] Ruscic, B.; Pinzon, R. E; von Laszewski, G.; Kodeboyina, D.; Burcat, A.; Leahy, D.; Montoy, D.; Wagner, A. F. Active thermochemical tables: thermochemistry for the 21st century *Journal of Physics: Conference Series* **2005**, *16*, 561.

[11] *IEEE Standard for Binary Floating-Point Arithmetic;* ANSI/IEEE Std 754-1985, 1985.

[12] *IEEE Standard for Floating-Point Arithmetic - Redline;* IEEE Std 754-2008 - Redline, 1985.

[13] Right-sizing precision: Unleashed computing: The need to right-size precision to save energy, bandwidth, storage, and electrical power. Gustafson, G. J. https://web.archive.org/web/20160606203112/http://www.johngustafson.net/presentations/Right-SizingPrecision1.pdf, **2015;** Presentation from web archive.

[14] Clark, M.A.; Babich, R.; Barros, K.; Brower, R.C.; Rebbi, C. Solving lattice qcd systems of equations using mixed precision solvers on gpus *Comp. Phys. Comm.* **2010**, *181*, 1517–1528.

[15] Grand, S. L.; Gtz, A. W.; Walker, R. C. SPFP: Speed without compromise — a mixed precision model for gpu accelerated molecular dynamics simulations *Comp. Phys. Comm.* **2013**, *184*, 374–380.

[16] Baboulin, M.; Buttari, A.; Dongarra, J.; Kurzak, J.; Langou, J.; Langou, J.; Luszczek, P.; Tomov, S. Accelerating scientific computations with mixed precision algorithms *Comp. Phys. Comm.* **2009**, *180*, 2526 – 2533; 40 YEARS OF CPC: A celebratory issue focused on quality software for high performance, grid and novel computing architectures.

[17] Asadchev, A.; Gordon, M. S. Mixed-precision evaluation of two-electron integrals by rys quadrature *Comp. Phys. Comm.* **2012**, *183*, 1563–1567.

[18] Rák, Á.; Cserey, G. The BRUSH algorithm for two-electron integrals on gpu *Chem. Phys. Lett.* **2015**, *622*, 92–98.

[19] Luehr, N.; Ufimtsev, I. S.; Martínez, T. J. Dynamic precision for electron repulsion integral evaluation on graphical processing units (gpus) *J. Chem. Theory Comput.* **2011**, *7*, 949–954.

[20] Ufimtsev, I. S.; Martínez, T. J. Quantum chemistry on graphical processing units. 1. strategies for two-electron integral evaluation *J. Chem. Theory Comput.* **2008**, *4*, 222–231.

[21] Ufimtsev, I. S.; Martínez, T. J. Quantum chemistry on graphical processing units. 2. direct self-consistent-field implementation *J. Chem. Theory Comput.* **2009**, *5*, 1004–1015.

[22] Ufimtsev, I. S.; Martínez, T. J. Quantum chemistry on graphical processing units. 3. analytical energy gradients, geometry optimization, and first principles molecular dynamics *J. Chem. Theory Comput.* **2009**, *5*, 2619–2628.

[23] Vogt, L.; Olivares-Amaya, R.; Kermes, S.; Shao, Y.; Amador-Bedolla, C.; Aspuru-Guzik, A. Accelerating resolution-of-the-identity second-order moller-plesset quantum chemistry calculations with graphical processing units *J. Phys. Chem. A* **2008**, *112*, 2049–2057.

[24] Vysotskiy, V. P.; Cederbaum, L. S. Accurate quantum chemistry in single precision arithmetic: Correlation energy *J. Chem. Theory Comput.* **2011**, *7*, 320–326.

[25] Knizia, G.; Li, W.; Simon, S.; Werner, H.-J. Determining the numerical stability of quantum chemistry algorithms *J. Chem. Theory Comput.* **2011**, *7*, 2387–2398.

[26] The metric for the amplitude convergence is the norm of the difference between the amplitudes from the successive iterations, i.e., $||T^{(i)} - T^{i-1}||$.

[27] The actual metric used in Molpro is the square of the norm of the difference between the amplitudes from the successive iterations, with the respective default threshold of $10^{-10}$.

[28] The metric for the amplitude convergence in GAMESS is the absolute value of the largest element of the difference between the amplitudes from the successive iterations, i.e., $\max|T^{(i)} - T^{i-1}|$.

[29] Epifanovsky, E.; Wormit, M.; Kuś, T.; Landau, A.; Zuev, D.; Khistyaev, K.; Manohar, P. U.; Kaliman, I.; Dreuw, A.; Krylov, A. I. New implementation of high-level correlated methods using a general block-tensor library for high-performance electronic structure calculations *J. Comput. Chem.* **2013**, *34*, 2293–2309.

[30] Kaliman, I.; Krylov, A. I. New algorithm for tensor contractions on multi-core CPUs, GPUs, and accelerators enables CCSD and EOM-CCSD calculations with over 1000 basis functions on a single compute node *J. Comput. Chem.* **2017**, *38*, 842–853.

[31] Krylov, A. I.; Gill, P. M. W. Q-Chem: An engine for innovation *WIREs Comput. Mol. Sci.* **2013**, *3*, 317–326.

[32] Y. Shao, Z. Gan, E. Epifanovsky, A.T.B. Gilbert, M. Wormit, J. Kussmann, A.W. Lange, A. Behn, J. Deng, X. Feng, D. Ghosh, M. Goldey, P.R. Horn, L.D. Jacobson, I. Kaliman, R.Z. Khaliullin, T. Kus, A. Landau, J. Liu, E.I. Proynov, Y.M. Rhee, R.M. Richard, M.A. Rohrdanz, R.P. Steele, E.J. Sundstrom, H.L. Woodcock III, P.M. Zimmerman, D. Zuev, B. Albrecht, E. Alguires, B. Austin, G.J.O. Beran, Y.A. Bernard, E. Berquist, K. Brandhorst, K.B. Bravaya, S.T. Brown, D. Casanova, C.-M. Chang, Y. Chen, S.H. Chien, K.D. Closser, D.L. Crittenden, M. Diedenhofen, R.A. DiStasio Jr., H. Do, A.D. Dutoi, R.G. Edgar, S. Fatehi, L. Fusti-Molnar, A. Ghysels, A. Golubeva-Zadorozhnaya, J. Gomes, M.W.D. Hanson-Heine, P.H.P. Harbach, A.W. Hauser, E.G. Hohenstein, Z.C. Holden, T.-C. Jagau, H. Ji, B. Kaduk, K. Khistyaev, J. Kim, J. Kim, R.A. King, P. Klunzinger, D. Kosenkov, T. Kowalczyk, C.M. Krauter, K.U. Laog, A. Laurent, K.V. Lawler, S.V. Levchenko, C.Y. Lin, F. Liu, E. Livshits, R.C. Lochan, A. Luenser, P. Manohar, S.F. Manzer, S.-P. Mao, N. Mardirossian, A.V. Marenich, S.A. Maurer, N.J. Mayhall, C.M. Oana, R. Olivares-Amaya, D.P. O'Neill, J.A. Parkhill, T.M. Perrine, R. Peverati, P.A. Pieniazek, A. Prociuk, D.R. Rehn, E. Rosta, N.J. Russ, N. Sergueev, S.M. Sharada, S. Sharmaa, D.W. Small, A. Sodt, T. Stein, D. Stuck, Y.-C. Su, A.J.W. Thom, T. Tsuchimochi, L. Vogt, O. Vydrov, T. Wang, M.A. Watson, J. Wenzel, A. White, C.F. Williams, V. Vanovschi, S. Yeganeh, S.R. Yost, Z.-Q. You, I.Y. Zhang, X. Zhang, Y. Zhou, B.R. Brooks, G.K.L. Chan, D.M. Chipman, C.J. Cramer, W.A. Goddard III, M.S. Gordon, W.J. Hehre, A. Klamt, H.F. Schaefer III, M.W. Schmidt, C.D. Sherrill, D.G. Truhlar, A. Warshel, X. Xu, A. Aspuru-Guzik, R. Baer, A.T. Bell, N.A. Besley, J.-D. Chai, A. Dreuw, B.D. Dunietz, T.R. Furlani, S.R. Gwaltney, C.-P. Hsu, Y. Jung, J. Kong, D.S. Lambrecht, W.Z. Liang, C. Ochsenfeld, V.A. Rassolov, L.V. Slipchenko, J.E. Subotnik, T. Van Voorhis, J.M. Herbert, A.I. Krylov, P.M.W. Gill, and M. Head-Gordon Advances in molecular quantum chemistry contained in the Q-Chem 4 program package *Mol. Phys.* **2015**, *113*, 184–215.

[33] Jagau, T.-C.; Bravaya, K. B.; Krylov, A. I. Extending quantum chemistry of bound states to electronic resonances *Annu. Rev. Phys. Chem.* **2017**, *68*, 525–553.

[34] Dreuw, A.; Wormit, M. The algebraic diagrammatic construction scheme for the polarization propagator for the calculation of excited states *WIREs Comput. Mol. Sci.* **2015**, *5*, 82–95.

[35] Ibrahim, K. Z.; Epifanovsky, E.; Williams, S.; Krylov, A. I. Cross-scale efficient tensor contractions for coupled cluster computations through multiple programming model backends *J. Parallel Distrib. Comput.* **2017**, *106*, 92–105.

[36] Levchenko, S. V.; Krylov, A. I. Equation-of-motion spin-flip coupled-cluster model with single and double substitutions: Theory and application to cyclobutadiene *J. Chem. Phys.* **2004**, *120*, 175–185.

[37] Curtiss, L. A.; Raghavachari, K.; Pople, J. A. Assessment of Gaussian-2 and density functional theories for the computation of enthalpies of formation *J. Chem. Phys.* **1997**, *106*, 1063–1079.

[38] The metric for the EOM amplitude convergence is the norm of the residual vector, $||\bar{H}R_k - \omega_k R_k||$, where $R_k$ and $\omega_k$ are EOM amplitudes and energies at a given iteration.